

# ENCODE Long Read RNA-Seq Analysis Protocol for Human Samples (v.1.0)

Prepared by Dana Wyman

June 2, 2019

Ali Mortazavi Lab, University of California, Irvine

## Contact Information

Dana Wyman  
2300 Biological Sciences III  
University of California Irvine  
Irvine, CA 92697-2300  
Telephone: (949) 824-8393  
Email: [dwyman@uci.edu](mailto:dwyman@uci.edu)

Ali Mortazavi  
2218 Biological Sciences III  
University of California Irvine  
Irvine, CA 92697-2300  
Telephone: (949) 824-6762  
Email: [ali.mortazavi@uci.edu](mailto:ali.mortazavi@uci.edu)

## I. Overview

The long-read RNA-seq data on the ENCODE portal was processed with the ENCODE DCC deployment of the TALON pipeline (available with documentation here: <https://github.com/ENCODE-DCC/long-read-rna-pipeline>).

This document describes 1) the steps used to generate the fastqs submitted to the DCC, and 2) individual tasks that make up the TALON pipeline. The software versions used can be found in Table 1.

The first two steps of the analysis pipeline, Circular Consensus (CCS), and pbtranscript Classify are part of the SMRTanalysis software suite available from Pacific Biosciences. They perform the following tasks, respectively:

1. Arrive at a consensus read of insert (ROI) sequence for each cDNA transcript
2. Remove primer/adaptor sequences and separate full-length ROIs from non full-length. Full-length ROIs are defined by the presence of a polyA tail and an adaptor sequence on each end.

The output of these steps are isoforms in the fastq/fasta format. These isoforms are mapped to the genome using Minimap2 to generate output in the sam/bam format. TranscriptClean is run on the mapped isoforms to correct remaining errors such as noncanonical splice junctions and microindels. TALON is run to annotate the transcripts and quantify their abundance.

Table 1: Referenced Software

| <b>Name</b>                       | <b>Version</b> | <b>Available from</b>   |
|-----------------------------------|----------------|---|
| DCC TALON pipeline implementation |                | <a href="https://github.com/ENCODE-DCC/long-read-rna-pipeline">https://github.com/ENCODE-DCC/long-read-rna-pipeline</a>                   |
| SMRTanalysis                      | 6.0.0          | Pacific Biosciences<br>( <a href="http://www.pacb.com/support/software-downloads/">http://www.pacb.com/support/software-downloads/</a> )  |
| Minimap2                          | 2.15           | <a href="https://github.com/lh3/minimap2">https://github.com/lh3/minimap2</a>   |
| TranscriptClean                   | 1.0.7          | <a href="https://github.com/dewyman/TranscriptClean">https://github.com/dewyman/TranscriptClean</a>                                       |
| TALON                             | 4.1            | <a href="https://github.com/dewyman/TALON">https://github.com/dewyman/TALON</a>   |
| Samtools                          | 1.3            | <a href="https://sourceforge.net/projects/samtools/files/samtools/1.3/">https://sourceforge.net/projects/samtools/files/samtools/1.3/</a> |
| fasta_to_fastq.pl                 | 3-24-2012      | <a href="https://github.com/ekg/fasta-to-fastq">https://github.com/ekg/fasta-to-fastq</a>   |

## II. Computational analysis

The programs used in this section are all part of the PacBio SMRTanalysis software suite.

### A. Obtaining reads of insert with Circular Consensus (CCS)

Multiple SMRT cells were sequenced per library in order to get > 1,000,000 raw reads, resulting in a **subreads.bam** file for each SMRT cell. CCS must be run on each of these files individually to generate consensus reads of insert (ROIs).

CCS is run using the following parameters:

```
ccs \  
  --noPolish \  
  --minLength=300 \  
  --minPasses=0 \  
  --minZScore=-999 \  
  --maxDropFraction=0.8 \  
  --minPredictedAccuracy=0.8 \  
  --minSnr=4 \  
  --reportFile ccs_report.txt \  
  subreads.bam \  
  ccs.bam
```

Each CCS run produces a bam file, **ccs.bam**. The sequences in these files are 'reads of insert' (ROIs), which represent the consensus sequence of each read.

### B. Isolating full-length, non-chimeric reads with pbtranscript Classify

The purpose of the pbtranscript Classify step is to identify full-length, non-chimeric (FLNC) reads based on the presence of a poly-A tail at the 3' end and adaptor sequences on each end. These reads will go on to form the basis for isoforms in the pbtranscript Cluster step. Non full-length (nfl) reads are filtered out at this point, but are not discarded. Instead, they are used to help with error correction in the pbtranscript Cluster step. A further role of pbtranscript Classify is to identify and filter out chimeric PacBio reads, which form when each end of a SMRTbell adaptor attaches to a different double-stranded cDNA molecule rather than to the blunt ends of the same one.

Pbtranscript Classify is run on each **ccs.bam** file separately using the following parameters:

```
pbtranscript.py classify \  
  ccs.bam  
  --min_seq_len 300 \  
  --cpus 16 \  
  --flnc SMRT_X_flnc.fasta  
  --nfl SMRT_X_nfl.fasta
```

The **SMRT\_X\_flnc.fasta** output file contains the full-length, non-chimeric ROIs, and **SMRT\_X\_nfl.fasta** contains the non full-length ROIs. Only the former is used in the downstream analysis. We convert **SMRT\_X\_flnc.fasta** to the fastq format using the following command:

```
perl fasta_to_fastq.pl SMRT_X_flnc.fasta > SMRT_X_flnc.fastq
```

The fastq files for all of the SMRT cells are concatenated together and submitted to the DCC . This file represents the starting input for their long read RNA-seq pipeline.

Note: the quality scores in this file are not meaningful.

## **C. Alignment to the reference genome**

We next align the FLNC reads to the GRCh38 XY human reference genome (<https://www.encodeproject.org/data-standards/reference-sequences/>). Run Minimap2 with the following parameters:

```
minimap2 -t 16 -ax splice -uf --secondary=no -C5 \  
  GRCh38.fa \  
  flnc.fastq \  
  > Aligned.out.sam
```

The output file will be called **Aligned.out.sam**:

## **D. Reference-based error correction**

### **I. Extract reference splice junctions**

TranscriptClean (next section) requires a file of reference splice junctions in order to correct noncanonical junctions in the PacBio Iso-seq isoforms. To get this file, we run a TranscriptClean utility on the GENCODE v29 comprehensive gene annotation

(reference chromosomes only). This file is available here:  
[https://www.gencodegenes.org/human/release\\_29.html](https://www.gencodegenes.org/human/release_29.html)

```
python ${TranscriptCleanPath}/accessory_scripts/get_SJs_from_gtf.py \  
  --f ../gencode.v29.annotation.gtf \  
  --g GRCh38.fa \  
  --o gencode_v29_SJs.tsv
```

The output file **gencode\_v29\_SJs.tsv**, contains splice junctions derived from the short reads. For each splice junction, it lists genomic location, strand, intron motif, and two additional placeholder columns (to match formatting to a type of STAR splice junction file).

## II. Error correction with TranscriptClean

Although the CCS process catches many of the errors found in PacBio transcripts, longer reads and/or those with fewer passes are still prone to mismatch and microindel errors. If these occur on the boundary of an intron, they may create the mistaken appearance of a novel splice junction. TranscriptClean is a Python program we developed to compare the sequences of mapped isoforms to the reference genome and correct likely errors. It can be downloaded from Github at <https://github.com/dewyman/TranscriptClean>. Run TranscriptClean version 1.0.7 on the FLNC reads using the parameters below. A file of reference splice junctions for (described in previous section) serves as a reference for correcting junctions. In addition, download a file of all common human variants from dbSNP Build150 (April 2017 release) in the VCF format to run TranscriptClean in variant-aware mode. This file can be found at: [https://www.ncbi.nlm.nih.gov/variation/docs/human\\_variation\\_vcf/](https://www.ncbi.nlm.nih.gov/variation/docs/human_variation_vcf/). When using variant-aware mode, mismatches and indels are compared to the VCF variant set, and are not modified in the event of a perfect variant match. Unmapped isoforms are discarded by TranscriptClean, and in the case of multi-mapping, only the primary mapping is used.

```
python TranscriptClean.py \  
  --sam Aligned.out.sam \  
  --genome GRCh38.fa \  
  --spliceJns gencode_v29_SJs.tsv \  
  --correctMismatches True \  
  --correctIndels True \  
  --variants 00-common_all.vcf.gz \  
  --maxLenIndel 5 \  

```

```
--maxSJOffset 5 \  
--outprefix libraryID
```

This command will generate two output files: **libraryID\_clean.sam** and **libraryID\_clean.fa**. These contain the reads with corrections made to remove microindels, mismatches, and noncanonical splice junctions as specified by the parameters. At this point, we remove any reads that still contain a noncanonical splice junction that is not in the reference annotation:

```
python filter_TC_output.py --f libraryID_clean.sam --o filtered_clean.sam
```

## **E. Annotate and quantify reads**

In order to compare long read platforms side by side and to track isoforms consistently across multiple datasets, we developed a technology-agnostic long read annotation tool called TALON. It is designed to annotate full-length reads as known or novel transcripts and also to report abundance for these transcripts. Corrected reads are passed into the TALON program, which is built around an SQLite database initialized to contain known genes, transcripts, and exon models from the GENCODE v29 GTF transcriptome annotation.

```
python initialize_talon_database.py \  
--f gencode.v29.annotation.gtf \  
--a gencode_v29 \  
--g mm10 \  
--l 300 \  
--idprefix ENCODEH \  
--5p 500 \  
--3p 300 \  
--o talon.db
```

In a TALON run, each input SAM transcript is compared to the existing transcript models in the database on the basis of its splice junctions, start, and end points. This allows us to not only assign a novel gene or transcript identity where appropriate, but to incorporate new transcript models in the TALON database while characterizing how they differ from known transcript models.

In order to take advantage of TALON's transcript filtering utilities, we run biological replicates through the program together. First, we create a comma-delimited config file (**config.csv**) with metadata about the samples. Column 1 is the unique dataset ID, column 2 is the cell type, column 3 is the platform, and column 4 is the input sam file. Here is an example:

```
Rep1,CellLine,PacBio-Sequel,Rep1/filtered_clean.sam  
Rep2,CellLine,PacBio-Sequel,Rep2/filtered_clean.sam
```

Next, we run TALON:

```
python talon.py --f config.csv \  
                --db talon.db \  
                --build hg38 \  
                --cov 0.9 \  
                --identity 0 \  
                --o cellLine
```

The TALON approach to quantification relies on each long read representing an individual transcript molecule sequenced, which allows us to quantify expression by simply counting the number of reads that were assigned to a particular transcript or gene and then converting these values to units of transcripts per million (TPM). To obtain a raw abundance file, run the following command on the TALON database:

```
python ${TALON}/post-TALON_tools/create_abundance_file_from_database.py \  
        --db talon.db \  
        -a gencode_v29 \  
        --build hg38 \  
        --o cellLine
```

The resulting file, **cellLine\_talon\_abundance.tsv**, can be used to compute gene-level expression values. For gene expression, we include all reads assigned to a locus except genomic transcripts, since even novel transcripts that did not meet the threshold to become a new transcript model are informative for the overall gene expression level. On the transcript level, however, we apply our TALON filters in order to avoid quantifying transcript models with insufficient evidence. Our filtering process uses the novelty labels assigned to each observed transcript model in order to remove likely artifacts. Observed transcripts that fully match counterparts in the GENCODE annotation are accepted immediately, but novel transcripts must be reproducibly

detected at least once across biological replicates in order to be included in the downstream analysis. Genomic transcripts are always removed, since they may represent genomic DNA contamination. To obtain a filtered abundance file for the datasets, run the following command:

```
python ${TALON}/post-TALON_tools/create_abundance_file_from_database.py \  
  --db talon.db \  
  -a gencode_v29 \  
  --build hg38 \  
  --filter \  
  -p pairings.csv \  
  --o cellLine
```

Here, pairings.csv is a file containing the dataset IDs (matching the config file) of the biological replicates. An example would be:

### **Rep1,Rep2**

It is also possible to generate a custom GTF transcriptome annotation for the samples from the TALON database. This file contains only transcript models from those samples that passed the TALON filters. The start and end coordinates used for these models are the ones first recorded in the database, so for known transcripts, that means the ones from GENCODE. To generate a filtered GTF, first run the filtering step to make a transcript whitelist:

```
python ${TALON}/post-TALON_tools/filter_talon_transcripts.py \  
  --db talon.db \  
  -a gencode_v29 \  
  -p pairings.csv \  
  --o whitelist.csv
```

Then, run the TALON GTF utility:

```
python ${TALON}/post-TALON_tools/create_GTF_from_database.py \  
  --db talon.db \  
  -a gencode_v29 \  
  -b hg38 \  
  --whitelist whitelist.csv \  
  --o cellLine_filtered
```

## **IV. References**

1. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
2. Gordon, S. P. et al. Widespread Polycistronic Transcripts in Fungi Revealed by Single-Molecule mRNA Sequencing. *PLoS ONE* 10, e0132628 (2015).
3. Li H.\*, Handsaker B.\*, Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R. and 1000 Genome Project Data Processing Subgroup (2009) The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics*, 25, 2078-9.
4. Wyman, D., TranscriptClean: A program for correcting mismatches, microindels, and noncanonical splice junctions in long reads, (2018), GitHub repository, <https://github.com/dewyman/TranscriptClean>